

Ontologies

Report

DBpedia

<Tuesday, 27 de July de 2010>

- Bruno Tavares; 31528

Index

Abstract.....	3
Introduction	3
Ontology.....	3
DBpedia.....	3
Geonames	4
Selected Ontologies	4
Methodology.....	4
DBpedia.....	4
GeoNames.....	5
Similarity	5
Examples	6
DBpedia.....	6
GeoNames.....	6
Similarity	6
Discussion	6
Appendix	8

Abstract

This work aims to gather geographic data from two different databases, Geonames and DBpedia, and calculate various similarities between the different elements, allowing an improvement of these databases and identify the overlap of these databases, allowing complement the information in each database.

Introduction

Ontology

Ontology is a schema to represent the knowledge based on a set of concepts, in a given field, and relations between them. There are several types of ontologies, such as disease (Unified Medical Language System), genes and gene products (Gene Ontology) and geographical (GeoNames) among others.

Ontologies usually share a similar structure; describing elements (instances) based on attributes (features, properties, etc.) and classes (concepts) and stipulate relationships between the various elements and classes. In a geographical ontology, one element could be Lisbon, which has as attributes a given population and geographic location (latitude and longitude), and belongs to a class such as “capital of a country”, and has relations like is-a “city”, which could be part-of “District of Lisbon”/“District of a Country” and in turn also be part-of “Portugal”/“Country”.

DBpedia

DBpedia is a project that aims to gather structured information from Wikipedia, providing access to this information in a Resource Description Framework (RDF) format, using a SQL-like query language (RDF query language) system called SPARQL (**SPARQL Protocol and RDF Query Language**). This database contains more than 3million entries, with almost half of these are in a consistent ontology, such as +300k people, +400k places or 140k organizations. For Portugal, this information is not very structured.

Geonames

GeoNames is a geographic ontology with over 8000000 entries, being categorized into nine classes (A [country, state, region,etc], P [city, village,etc], etc), and nearly 650 subclasses (independent political entity, populated place, etc). Like other ontologies, Geonames has attributes on each element (District of Lisbon), such as its population and geographical coordinates. And the relationships between different elements, in the case of Portugal these relationships are poorly annotated.

Selected Ontologies

The selected ontologies to perform this work were Geonames (reliable source of geographic information) and DBpedia (source constantly updated).

Methodology

DBpedia

Because of the weak structure of the data from DBpedia, the process of collecting such data was implemented in two phases. In the first phase was created the data structure based on the database dump from the Portuguese version of Wikipedia. In the second phase were executed SPARQL queries to DBpedia to gather information.

In the first phase, were executed these steps:

- List the districts of Portugal;
- Collect the id's of the links from the pages of each district;
- Replace the id's for page titles;
- List all id's and titles for the pages with the same title page listed above (One id link could correspond to a not geographic type, but different entries could have the same title and one of this correspond to a geographic type);

- Insert in the previous list id's and titles of pages whose title is similar to those previously collected but with a termination "`_(<term>)`" (Some pages have the title `<municipality>_(<district>)` but in database that link appears with the title `<municipality>`);
- Filter links by type ("`Municípios_de_Portugal`");
- Remove the duplicate and incorrect entries (Some pages are the type "`Municípios_de_Portugal`" but not are a geographical location, for example, it could be an entry for "`clube_desportivo_<name>`");
- Inserts new entries that were not previously entered (for some reason the municipalities of Vila Real belongs to the type "`Distrito_de_Vila_Real`" unlike all the others municipalities that belong to type "`Municípios_de_Portugal`", it should be noted that, exception for these municipalities, are rare cases of municipalities that belong to the type "`Distrito_de_<distrito>`");
- List the parishes from each municipality;

Based on this structure is then possible to collect information about each geographical place. These data are collected in a cyclical manner, using queries-SPARQL. Two queries are made to DBpedia, the first is the redirect of the pages title from Portuguese to English, this process needs to be done because the structure was made from the Portuguese version of Wikipedia and so the pages title are in Portuguese and secondly the DBpedia do not have a Portuguese version. The other query can thus collect information from locations such as population, latitude and longitude.

GeoNames

The process of collecting information from Geonames is facilitated by the Web Services available. Two distinct Web Services are used, the first (Search) returns all entries in accordance with certain attributes (type and country). The other Web Service (hierarchy) returns the parent of an entry.

Similarity

With the collection of information from these two databases, is now permitted to pairing them, thus enabling a comparison of different geographic ontologies. Was done three different similarities measure, a comparison of the names of localities using the Levenshtein measure, a difference of populations and geographical distance.

Examples

DBpedia

```
SELECT * WHERE {  
  
<http://dbpedia.org/resource/<page_title>> ?type ?value .  
  
Filter( ?type = dbpprop:populationTotal ||  
  
?type = geo:lat ||  
  
?type = geo:long ) }
```

GeoNames

```
http://ws.geonames.org/search?country=PT &style=full&featureCode=<type>  
  
http://ws.geonames.org/hierarchy?geonameId=<geonameId>
```

Similarity

```
delta_population = abs( log( pop1 / pop2 ) )  
  
distance = 6371 * acos( ( ( cos(lon1 - lon2 ) * ( cos(lat1 - lat2 ) + cos(lat1 + lat2 ) ) ) + ( cos(lat1 - lat2 ) -  
cos(lat1 + lat2 ) ) ) / 2 )
```

Discussion

The biggest critic I can do to my work is that ontologies similarities techniques were not applied, techniques that weigh the similarity measures of a pair, with the similarities between the parents of this pair, but the process executed was done in order to allow these techniques in the future. Regarding the process executed to structure the data from DBpedia, to emphasize the fact that part of the dump of wikipedia is not consistent, eg two different pages (id's) can have the same title, this should mean that one URI could match different pages, what does not. Moreover, there are links hardly annotated, for

example in the Webpage of "Distrito_de_Vila_Real" there is a link to the municipality of "Mesão_Frio" and in the database the title of this page is "Mesão_Frio_ (Vila_Real).

One of the positive points of this work is to enable the improve of the structures of both databases, on one hand the GeoNames has some gaps in their structure between the administrative divisions and geographic locations, on the other hand DBpedia structure is virtually nonexistent and the executed work was curated by hand.

Speaking now on the results of similarity, I have more than 113 million results, being that, for each pair (more than 72M pairs) there is at least the result of the Levenshtein measure. An example of the results of similarity can be "Caldas da Rainha," because this location has all the information in both databases, containing the following results: for the method Levenshtein the similarity is equal to 1, the names are equal; in relation to geographical distance the obtained value was 0.8, this means that the distance between this location in the different databases is only 0.8 km; for differences in population the result was 0.6, indicating that the populations have the same magnitude but slightly different (26k for Geonames and 48k for DBpedia).

As mentioned previously, these results could allow a complement of data in both databases. One of the points where this can be well executed will be on the types of locations, the good structure created for DBpedia, can be extrapolated to the Geonames. In the example given above we can then indicate that "Caldas da Rainha" is a municipality, for this is necessary to make a match of the types in Geonames. Giving a concrete example, in the districts we have similarities values of 1 for the method Levenshtein, exception of two locations that I will refer later on, so you may have a correlation between the type, in DBpedia the types are "Distrito" while in Geonames are "first-order administrative division". From this correlation, and because we do not have all the information in both databases, we can still connect the "Distrito de Lisboa" (DBpedia) with the "Lisbon" (Geonames), Levenshtein similarity of 0.2 and the "Distrito de Bragança" (DBpedia) with "Bragança" (Geonames), Levenshtein similarity of 0.4.

Appendix

```
from MySQLdb import connect, IntegrityError
from math import atan, cos, log, radians
from xml.parsers.expat import ExpatError
from socket import setdefaulttimeout
from urllib import urlopen, quote
from xml.dom.minidom import parse
from configobj import ConfigObj
from time import sleep
from os import system
from sys import argv

setdefaulttimeout( 0.5 * 60 )

def insertSQL( query, database = db ):
    conn = connect( host = "localhost", user = user, passwd = password, db = database )
    new_cursor = conn.cursor()
    new_cursor.execute( query )
    new_cursor.close()
    conn.commit()
    conn.close()

def querySQL( query, database = db ):
    conn = connect( host = "localhost", user = user, passwd = password, db = database )
    new_cursor = conn.cursor()
    new_cursor.execute( query )
    result = new_cursor.fetchall()
    new_cursor.close()
    conn.commit()
    conn.close()
    return result

def api( url ):
    f = urlopen( url )
    if doc = parse( f ):
        return ExpatError: doc = api( url )
    f.close()
    except IOError: doc = api( url )
    return doc

def createDB( database = db ):
    conn = connect( host = "localhost", user = user, passwd = password, db = db )
    new_cursor = conn.cursor()
    new_cursor.execute( "drop database if exists " + database )
    new_cursor.execute( "create database " + database )
    new_cursor.close()
    conn.commit()
    conn.close()

def createTables():
    conn = connect( host = "localhost", user = user, passwd = password, db = db )
    new_cursor = conn.cursor()
    new_cursor.execute( "drop table if exists Deoname" )
    new_cursor.execute( "drop table if exists Acc_Deoname" )
    new_cursor.execute( "drop table if exists GeoDeoname" )
    new_cursor.execute( "drop table if exists Acc_GeoDeoname" )
    new_cursor.execute( "drop table if exists Updeals" )
    new_cursor.execute( "drop table if exists Pair_gsm_pt" )
    new_cursor.execute( "drop table if exists Pair_gsm_db" )
    new_cursor.execute( "drop table if exists Similarity_gsm_db" )
    new_cursor.execute( "drop table if exists Similarity_pt_db" )
    new_cursor.execute( "create table Deoname (id int primary key auto_increment, name varchar(100), lat float, lng float, pop int, type varchar(80), geoNameId int unique, index(name(4)), index(lat), index(lng), index(type))" )
    new_cursor.execute( "create table Acc_Deoname (id int primary key auto_increment, name varchar(100), lat float, lng float, pop int, type varchar(80), geoNameId varchar(100) unique, index(name(4)), index(lat), index(lng), index(type))" )
    new_cursor.execute( "create table GeoDeoname (id int primary key auto_increment, term1 int, term2 int, distance int, index(term1))" )
    new_cursor.execute( "create table Updeals (id int primary key auto_increment, name varchar(100), lat float, lng float, pop int, type varchar(80), dbDealId varchar(100), index(name(4)), index(lat), index(lng), index(type))" )
    new_cursor.execute( "create table Acc_Updeals (id int primary key auto_increment, term1 int, term2 int, distance int, index(term1))" )
    new_cursor.execute( "create table Similarity_gsm_pt (id int primary key auto_increment, pairId int, method varchar(32), similarity float, relation varchar(32), index(pairId), index(similarity), index(method))" )
    new_cursor.execute( "create table Pair_gsm_db (id int primary key auto_increment, term1 int, term2 int, index(term1), index(term2))" )
    new_cursor.execute( "create table Similarity_gsm_db (id int primary key auto_increment, term1 int, term2 int, method varchar(32), similarity float, relation varchar(32), index(pairId), index(similarity), index(method))" )
    new_cursor.close()
    conn.commit()
    conn.close()

def levenshtein( s1, s2 ):
    l1 = len( s1 )
    l2 = len( s2 )
    matrix = [ range( l1 + 1 ) ] * ( l2 + 1 )
    for xx in range( l2 + 1 ):
        matrix[ xx ] = range( xx, xx + l1 + 1 )
    for xx in range( 0, l2 ):
        for xx in range( 0, l1 ):
            if s1[ xx ] == s2[ xx ]:
                matrix[ xx + 1 ][ xx + 1 ] = min( matrix[ xx + 1 ][ xx ] + 1, matrix[ xx ][ xx + 1 ] + 1, matrix[ xx ][ xx ] )
            else:
                matrix[ xx + 1 ][ xx + 1 ] = min( matrix[ xx + 1 ][ xx ] + 1, matrix[ xx ][ xx + 1 ] + 1, matrix[ xx ][ xx ] + 1 )
    return l1 - float( matrix[ l2 ][ l1 ] ) / max( l1, l2 )

def delta_geo( geo1, geo2 ):
    """Return the difference between the two latitudes or longitudes.
    """
    return abs( geo1 - geo2 )

earth_radius = 6371
def distance( lat1, lon1, lat2, lon2 ):
    """Return the minimum distance in Earth's spherical surface between the points
    defined by coordinates (lat1, lon1) and (lat2, lon2). The values must be given
    in degrees.
    """
    # Convert degrees to radians
    # = x: latitude
    # = y: longitude
    x1, y1, x2, y2 = ( radians( i ) for i in ( lat1, lon1, lat2, lon2 ) )
    p1 = cos( y1 - y2 )
    p2 = cos( x1 - x2 )
    p3 = cos( x1 + x2 )
    inner_product = ( ( p1 + ( p2 + p3 ) ) + ( p2 - p3 ) ) / 2
    return earth_radius * acos( inner_product )
```


[illegible]

[illegible]

[illegible]